# A taste of
# Supercomputer Reliability Research
# at Sandia

**Red Storm Quarterly Review**
**Oct 25, 2007**

**Jon Stearley**
*jrstear@sandia.gov*
**Scalable Systems Architecture (1422)**

# Reliability is hard!

**System Facts:**

- 124.42 teraOPS theoretical peak performance
- **12,960** compute nodes,
  320 + 320 service and I/O nodes
- 40 terabytes of DDR memory
- 340 terabytes of disk storage
- Linux/Catamount Operating Systems
- Approximately 3500 ft2 including disk systems
- <2.5 megawatts of power and cooling
- **3,710** Linux computers used to control this beast
- 14,240 high-speed network interfaces

Hmmm, **which wire…** is loose?

**Mission Critical**

**+**

**Many points of failure**

**+**

**Complex and dynamic interdependencies**

**=**

**Rich research area!**

Sandia
National
Laboratories

# System Logs

Are:    Ubiquitous!  Informational! Vast!



How do you find the few lines of key information among thousands of log files and millions of lines of time-stamped text???
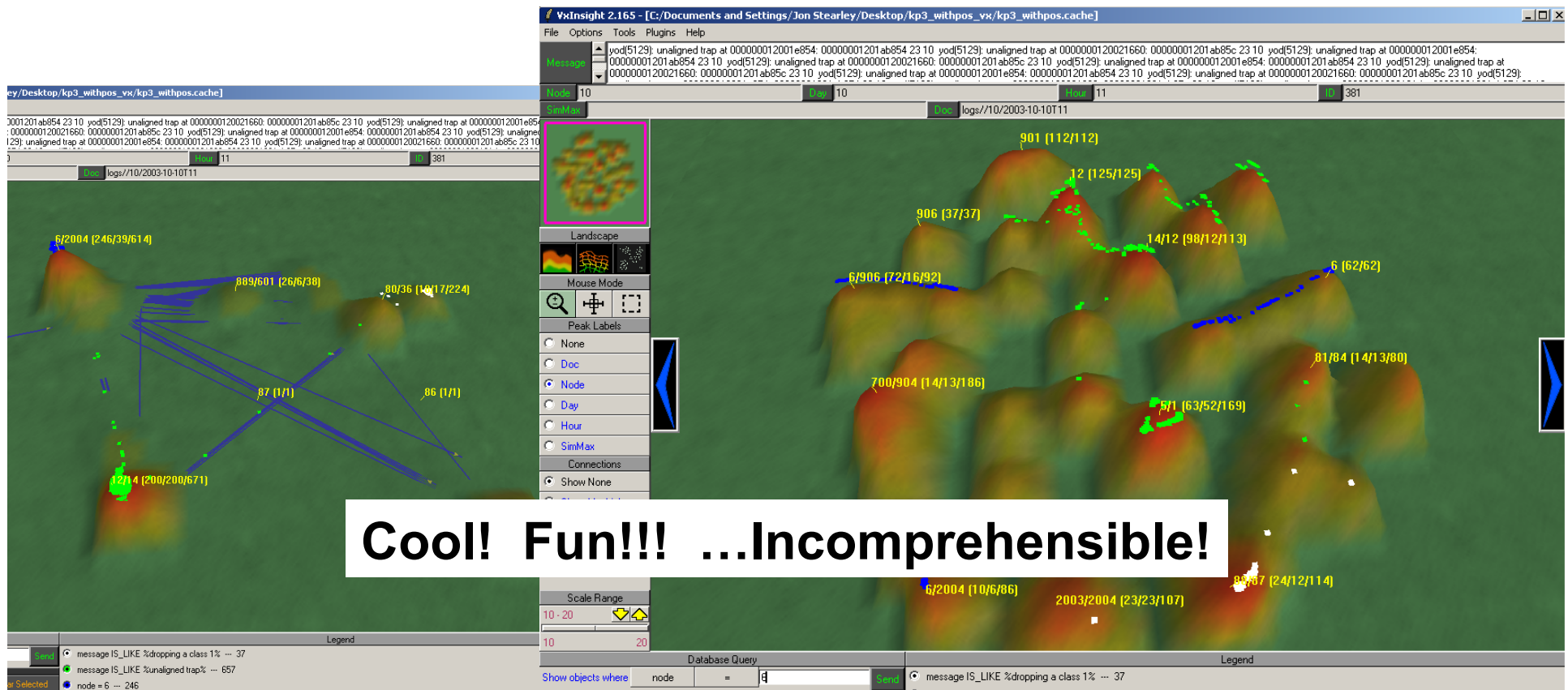
**Key Idea:**

Similar computers…
correctly performing similar work…
should produce similar logs.

(Anomalies warrant investigation.)

# Latent Semantic Analysis

1. Calculate logfile-logfile similarities (via SVD)
2. Cluster (VxOrd)
3. Explore themes (VxInsight/Threatview)
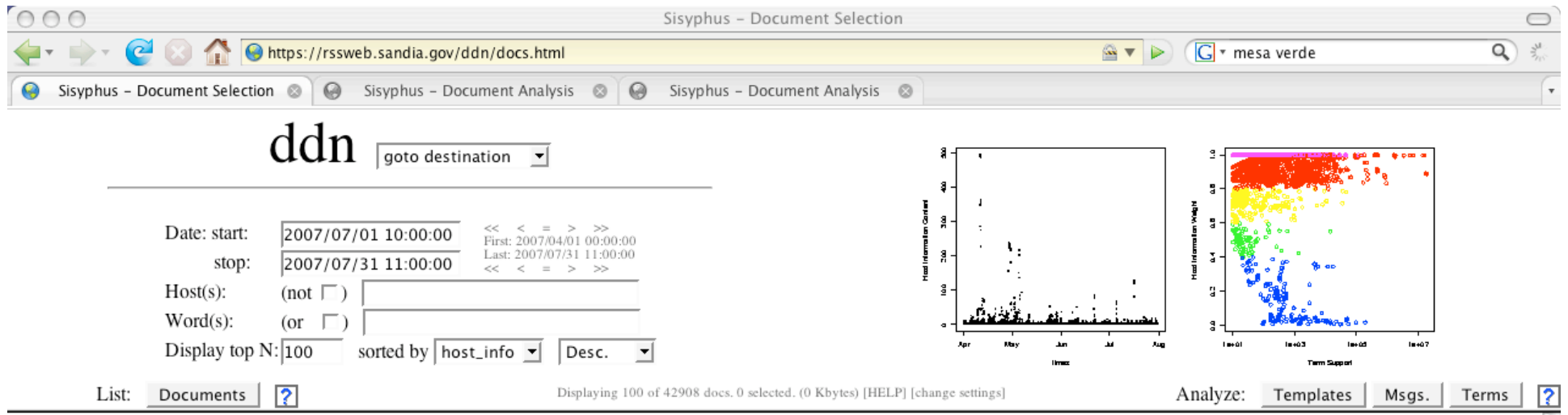


**Cool!  Fun!!!  …Incomprehensible!**

# Finding Needles in a Craystack

1. Which files contain <u>useful</u> information?
2. Which words convey <u>useful</u> information?
3. Any patterns?

**To be Useful,**

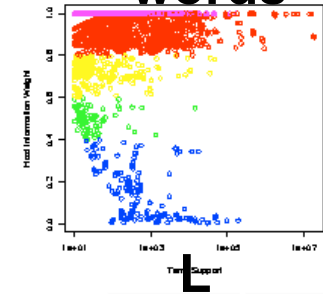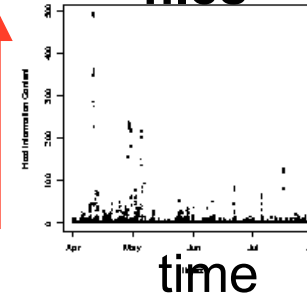**It Must be Understandable**

**(to the sysadmins)**

# 1. Which files contain <u>useful</u> information?

Sisyphus – Document Selection

https://rssweb.sandia.gov/ddn/docs.html · mesa verde

ddn   goto destination

**files**   **words**

$|(GL)_j|$   **G**

time   **L**

Date: start: 2007/07/01 10:00:00    << < = > >>   First: 2007/04/01 00:00:00
stop: 2007/07/31 11:00:00    Last: 2007/07/31 11:00:00    << < = > >>
Host(s): (not ☐)
Word(s): (or ☐)
Display top N: 100   sorted by host_info ▾ Desc. ▾

List: Documents ?   Displaying 100 of 42908 docs. 0 selected. (0 Kbytes) [HELP] [change settings]   Analyze: Templates Msgs. Terms ?

| | YYYY/MM/DD/HH HOST | bytes | lines | host_info | time_info | doc_info |
|---|---|---|---|---|---|---|
| ☐ | docs/2007/07/16/10/10.1.0.49 | 1842113 | 11251 | 130.493 | 102.764 | 111.610 |
| ☐ | docs/2007/07/16/09/10.1.0.49 | 1867390 | 11437 | 129.803 | 102.591 | 111.241 |
| ☐ | docs/2007/07/16/11/10.1.0.49 | 1816549 | 11125 | 129.538 | 102.359 | 111.008 |
| ☐ | docs/2007/07/16/12/10.1.0.49 | 1704339 | 10437 | 126.824 | 100.048 | 108.612 |
| ☐ | docs/2007/07/16/08/10.1.0.49 | 1481224 | 9068 | 120.769 | 95.549 | 103.661 |
| ☐ | docs/2007/07/16/13/10.1.0.49 | 430320 | | | | |
| ☐ | docs/2007/07/05/09/10.1.0.12 | 288005 | | | | |
| ☐ | docs/2007/07/05/09/10.1.0.16 | 158502 | | | | |
| ☐ | docs/2007/07/05/09/10.1.0.11 | 77539 | | | | |
| ☐ | docs/2007/07/05/09/10.1.0.6 | 17907 | | | | |
| ☐ | docs/2007/07/22/14/10.1.0.35 | 4430 | | | | |
| ☐ | docs/2007/07/22/05/10.1.0.35 | 4210 | | | | |
| ☐ | docs/2007/07/27/12/10.1.0.47 | 4887 | | | | |
| ☐ | docs/2007/07/22/13/10.1.0.35 | 3324 | | | | |
| ☐ | docs/2007/07/05/10/10.1.0.28 | 11386 | | | | |
| ☐ | docs/2007/07/29/19/10.1.0.47 | 3746 | | | | |
| ☐ | docs/2007/07/05/12/10.1.0.38 | 11966 | | | | |
| ☐ | docs/2007/07/27/14/10.1.0.47 | 3526 | | | | |
| ☐ | docs/2007/07/23/17/10.1.0.35 | 2660 | | | | |
| ☐ | docs/2007/07/05/12/10.1.0.46 | 10482 | | | | |
| ☐ | docs/2007/07/05/12/10.1.0.51 | 10653 | | | | |
| ☐ | docs/2007/07/05/09/10.1.0.7 | 9918 | | | | |
| ☐ | docs/2007/07/05/14/10.1.0.46 | 9172 | | | | |
| ☐ | docs/2007/07/05/14/10.1.0.38 | 9637 | | | | |
| ☐ | docs/2007/07/05/14/10.1.0.51 | 9480 | | | | |
| ☐ | docs/2007/07/24/10/10.1.0.35 | 2218 | | | | |
| ☐ | docs/2007/07/05/13/10.1.0.46 | 6495 | 68 | 15.809 | 16.447 | 15.882 |
| ☐ | docs/2007/07/05/12/10.1.0.7 | 6688 | 72 | 15.857 | 16.686 | 15.964 |
| ☐ | docs/2007/07/05/13/10.1.0.28 | 6440 | 68 | 15.833 | 16.230 | 15.729 |
| ☐ | docs/2007/07/05/12/10.1.0.28 | 6130 | 64 | 15.782 | 16.020 | 15.688 |
| ☐ | docs/2007/07/05/13/10.1.0.38 | 6800 | 72 | 15.746 | 16.551 | 15.868 |

abnormal

"interestingness"

normal

"interestingness"
(aka "information") is *purely*
mathematical (=$|(GL)_j|$).

$G_{i,j}=1+H_i$  ,  $L=\log_2(tf_{i,j})$
$H_i=\sum_j p_{ij}\log_2(p_{ij})/\log_2(d)$
  where $p_{ij}= tf_{i,j}/\sum_j tf_{i,j}$
  and $tf_{i,j}$ is how many times the
    i'th word occurs in the
    j'th file

Done   rssweb.sandia.gov

How do you find the few lines of key information among thousands of log files and millions of lines of time-stamped text???

2. Which words convey <u>useful</u> information?

<docs/2007/07/16/10/10.1.0.49>   **Email URL**    8 templates, minsup=20. 5274 terms. 0 selected. [change settings]    Analyze: | Templates | Msgs. | Terms

```
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Medium Error Disk 4G 3KT1HVCG Key: 3 ASC 16 ASCQ 0 FRU D2 Sense 80008E Info 0889A800
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889a800 LUN 7, 00000090 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889aa00 LUN 7, 00000091 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889ac00 LUN 7, 00000092 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889ae00 LUN 7, 00000093 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b000 LUN 7, 00000094 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b200 LUN 7, 00000095 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b400 LUN 7, 00000096 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b600 LUN 7, 00000097 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000090 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000091 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000092 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000093 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000094 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000095 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000096 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, ...0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Medium Error Disk 4G 3K...               ...nse 80008E Info 02244600
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G add...            ...0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G add... 2244800 LUN 6, 00011224 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2244a00 LUN 6, 00011225 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2244c00 LUN 6, 00011226 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2244e00 LUN 6, 00011227 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245000 LUN 6, 00011228 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245200 LUN 6, 00011229 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245400 LUN 6, 0001122a DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011223 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011224 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011225 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011226 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011227 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011228 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011229 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:05 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122a DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Medium Error Disk 4G 3KT1HVCG Key: 3 ASC 16 ASCQ 0 FRU D2 Sense 80008E Info 02245600
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245600 LUN 6, 0001122b DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245800 LUN 6, 0001122c DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245a00 LUN 6, 0001122d DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245c00 LUN 6, 0001122e DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245e00 LUN 6, 0001122f DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246000 LUN 6, 00011230 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246200 LUN 6, 00011231 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246400 LUN 6, 00011232 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122b DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122c DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0...
Jul 16 10:00:07 10.1.0...
Jul 16 10:00:07 10.1.0...
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011231 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011232 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:17 10.1.0.49 local7 info INT_DG Medium Error Disk 4G 3KT1HVCG Key: 3 ASC 16 ASCQ 0 FRU D2 Sense 80008E Info 02246600
Jul 16 10:00:17 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246600 LUN 6, 00011233 DLR:0, DLG:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
```

**A gold mine!!!**

**2. Which words convey <u>useful</u> information?**

**2. Which words convey <u>useful</u> information?**

# 2. Which words convey <u>useful</u> information?

# 3. Are there any patterns?

**Time patterns.**

**Linewise word patterns via sequencing (Teiresias), clustering (SLCT), and association (Apriori).**

| | ID | count | median | stddev | regexp. |
|---|---|---|---|---|---|
| | 0 | 113 | 0 | 57 | OUTLIERS |
| | 1 | 40 | 15 | 100 | daemon info llrd 5640 : llrd: nid00192 - - 17/Oct/2007 * "POST /RPC2 HTTP/1.0" 200 - |
| | 2 | 20 | 0 | 389 | kern * kernel: * * slow * * |
| | 3 | 9 | 301 | 472 | kern err kernel: LustreError: * * * * |
| | 4 | 47 | 13 | 102 | kern alert kernel: LustreError: dumping log to * |
| | 6 | 4 | 760 | 853 | kern * kernel: * dumping log to * |
| | 7 | 6 | 602 | 16 | kern * kernel: * * * * * |
| | 8 | 86 | 22 | 132 | kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000 |
| | 18 | 86 | 22 | 132 | kern warning kernel: end_request: I/O error, dev sde, sector * |
| | 20 | 50 | 0 | 97 | kern warning kernel: Call Trace:{schedule_timeout+243} {process_timeout+0} |
| | 21 | 36 | 0 | 79 | kern warning kernel: Call Trace:{:libcfs:libcfs_nid2str+178} {:ost:ost_brw_write+2000} |
| | 22 | 2 | 301 | 0 | kern warning kernel: Call Trace:{:libcfs:libcfs_nid2str+178} * |
| | 23 | 2 | 600 | 0 | kern warning kernel: Call * {:ost:ost_brw_write+2000} |

Oct 17 05:04:06 nid00187 kern crit kernel: LDISKFS-fs error (device sde2) in ldiskfs_setattr: Readonly filesystem
Oct 17 05:04:12 nid00187 kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000
Oct 17 05:04:12 nid00187 kern warning kernel: end_request: I/O error, dev sde, sector 778694416
Oct 17 05:04:12 nid00187 kern err kernel: Buffer I/O error on device sde2, logical block 7372802
Oct 17 05:04:12 nid00187 kern warning kernel: lost page write due to I/O error on sde2
Oct 17 05:04:12 nid00187 kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000
Oct 17 05:04:12 nid00187 kern warning kernel: end_request: I/O error, dev sde, sector 779218704
Oct 17 05:04:12 nid00187 kern err kernel: Buffer I/O error on device sde2, logical block 7438338
Oct 17 05:04:12 nid00187 kern warning kernel: lost page write due to I/O error on sde2
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(lustre_fsfilt.h:255:fsfilt_commit_wait()) slow journal start 51s
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(filter_io_26.c:707:filter_commitrw_write()) slow commitrw commit 3511s
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(filter_io_26.c:707:filter_commitrw_write()) previously skipped 5 similar messages
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(service.c:583:ptlrpc_server_handle_request()) request 527 opc 4 from U3-1251@ptl processed in 3511s trans 0 rc -5/-5
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(service.c:583:ptlrpc_server_handle_request()) previously skipped 7 similar messages
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(watchdog.c:320:lcw_update_time()) Expired watchdog for pid 6388 disabled after 3511.0309s
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6339:0:(watchdog.c:320:lcw_update_time()) Expired watchdog for pid 6339 disabled after 3511.4820s
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(watchdog.c:320:lcw_update_time()) previously skipped 7 similar messages

# Production Impacts

**Sisyphus has found a wide range of problems:**

**Failures:**
Disks and controllers
Network interfaces
Power supplies
Memory

**Misconfigurations:**
Performance-decreasing BIOS setting
Overhead-increasing RAID controller setting
Inconsistent software versions across nodes
Faulty software configuration

**Problematic user behavior:**
Unbalanced disk RAID stripe usage
Inappropriate remote monitoring

**Which has enabled focused reactive and proactive responses.**

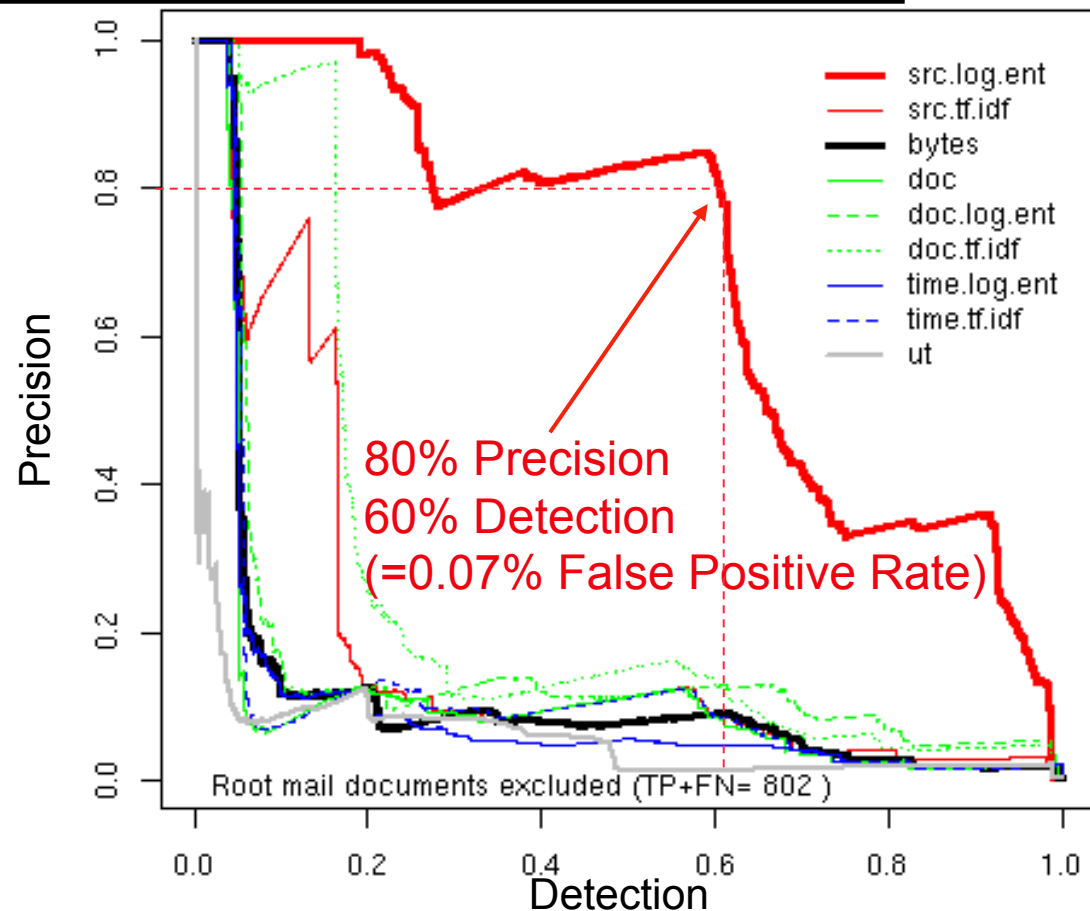*See http://www.cs.sandia.gov/sisyphus for more info.*

## 33 Unsupervised Classifiers Tested!

NWCC/Spirit Data
512 Nodes, 23 Days
8.3M log messages
36K terms, 243K docs
3.9K emails!
P=62 ; 802   N=243K



80% Precision
60% Detection
(=0.07% False Positive Rate)

Legend:
- src.log.ent
- src.tf.idf
- bytes
- doc
- doc.log.ent
- doc.tf.idf
- time.log.ent
- time.tf.idf
- ut

Root mail documents excluded (TP+FN= 802 )

Precision (y-axis)
Detection (x-axis)

True Class:

|              | P   | N   |
|--------------|-----|-----|
| Alarm Class: P | TP  | FP  |
| Alarm Class: N | FN  | TN  |

**TP**=True Positives
**FP**=False Positives
**FN**=False Negatives
**TN**=True Negatives

**Metrics:**
Alarm Precision = TP/(TP+FP)
Event Detection = TP/(TP+FN)

# RAS Metrics

RAS = Reliability, Availability, Serviceability

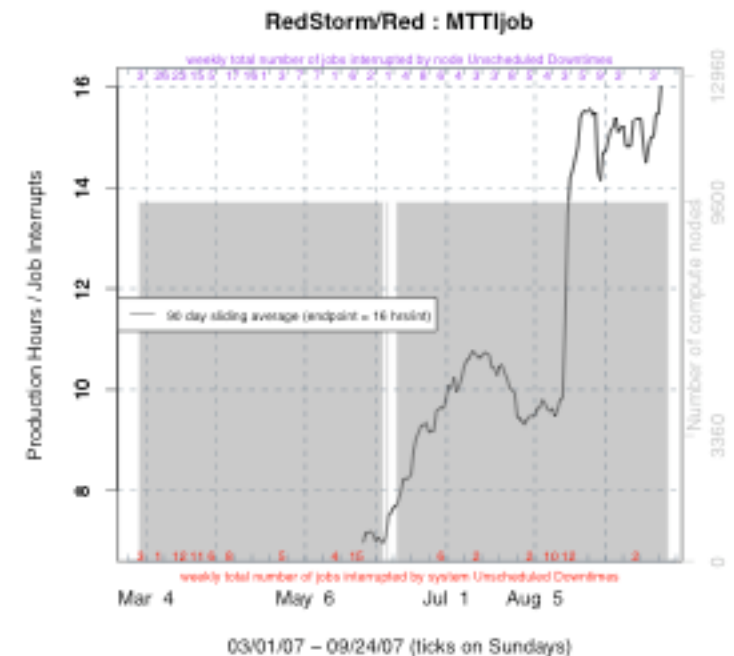**Good science is measurable!**
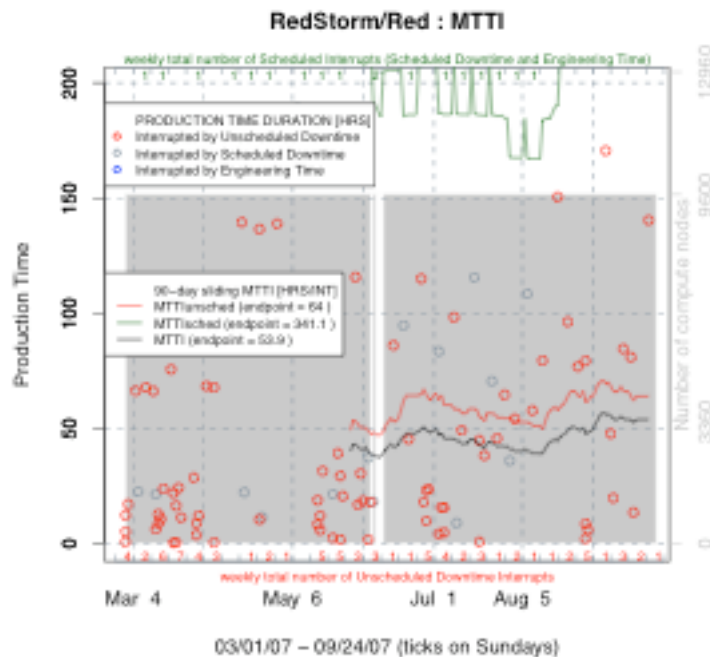
**But…**

**NO STANDARD is currently used
for measuring supercomputer RAS!!!**

# RAS Metrics: Challenges

**Difficult to agree on:**

- How to define failure (or interrupt).

- How to measure reliability.

- The need and method to change the processes and procedures involved.

# RAS Metrics: Plans

The Tri-laboratory Linux Compute Cluster presents a fantastic RAS metrics opportunity:

- Same vendor
- Same hardware
- Same system software

Our FY08 goal is to produce a specification and reference implementation for TLCC RAS metrics.
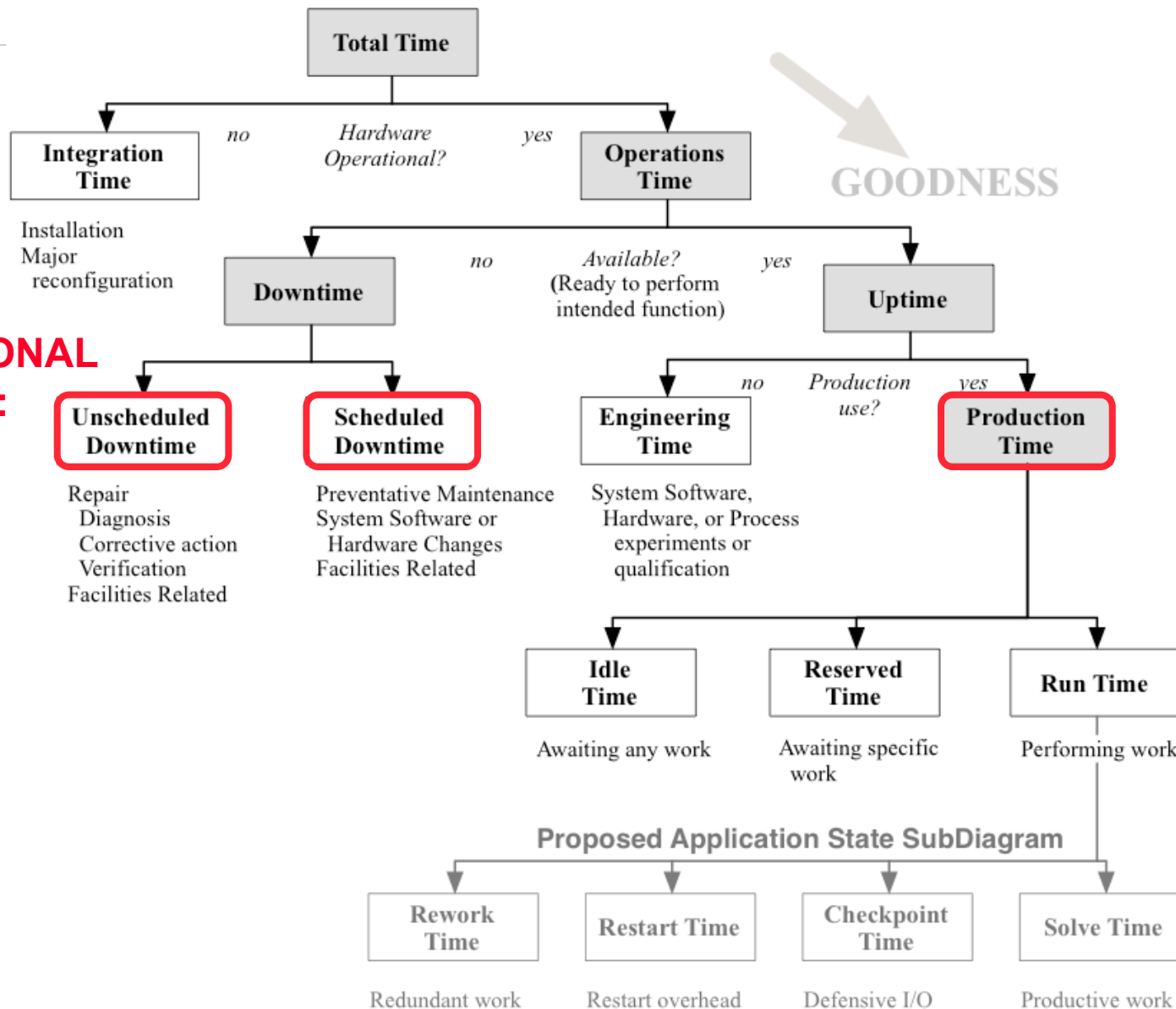
Key Idea:

track **OPERATIONAL CONTEXT**
for every node, at all times

**Tri-lab-developed Component State Diagram** *(Based on SEMI-E10)*
Each component is in exactly one non-grey state at all times.

# Summary

**Supercomputer RAS is a rich research area.**

**Sandia is making significant contributions.**

Applications, Operating systems, I/O systems, System
architectures, Device control, Networking, Metrics,
and Detection and prediction on logs and real-valued data.

(I have given you only a taste)

**Standard RAS metrics are essential.**

For improved RAS research, engineering, and operation.